# Optimal admission, routing and service assignment control: the case of single buffer queues[*]

Eitan ALTMAN and Bruno GAUJAL
INRIA
2004 Route des Lucioles
BP93, 06902 Sophia-Antipolis Cedex
France
Tel. +33 4 92 38 77 86 / 76 36
Fax +33 4 92 38 78 58
Email {altman,gaujal}@sophia.inria.fr

Arie HORDIJK
Leiden University
Dept. of Mathematics and Computer Science
P.O.Box 9512, 2300 RA Leiden
The Netherlands
Tel. +31 71 5277146
Fax. +31 71 5276985
Email hordijk@wi.leidenuniv.nl

Ger KOOLE
Vrije Universiteit
Faculty of Mathematics and Computer Science
De Boelelaan 1081a, 1081 HV Amsterdam
The Netherlands
Tel. +31 20 444 7755
Fax +31 20 444 7653
Email koole@cs.vu.nl

## Abstract

We consider the problem of optimal routing of arriving packets into $N$ servers having no waiting room. Packets that are routed to a busy server are lost. We consider two problems where the objective is to maximize the expected throughput (or equivalently, minimize the

---

loss rate). We assume that the controller has no information on the state of the server. We establish the optimality of the so called "balanced" policies, for exponential service times and general stationary arrival processes, which include, in particular, the interrupted Poisson process, Markov modulated Poisson Process (MMPP) and Markov arrival process (MAP). Based on this solution, we solve the dual problem of optimal assignment of a single server to several single server queues to which packets arrive according to Poisson processes. This general model is then applied to solve an optimal scheduling problem for robots of Web search engines.

# 1   Introduction

In a recent paper [8] Koole used the theory of MDPs with partial information to solve the problem of optimal routing packets into $N$ servers with different service time. The objective was to minimize losses or maximize the throughput, and policies are restricted to those where no information is available on the state of the servers. For the special case of symmetrical servers and arrival processes, the optimality of the round-robin policy was established. For the case of 2 servers, it was shown that a periodic policy whose period has the form $(1, 2, 2, 2, \ldots, 2)$ is optimal, where 2 means sending a packet to the faster server. Similar results have been obtained in the dual model of server assignment in [5]. The results in [5, 8] heavily depend on the Markovian structure of the arrivals.

In this paper we use a novel approach developed in [1] (and applied to other control problems in [2, 3]) that enables us to relax the Markovian assumptions (more precisely, the assumption that the inter-arrival times are i.i.d.). We begin by considering two types of problems: an admission into a single server, and the routing problem studied in [8]. We establish the optimality of the so called "balanced" policies [6] (see also [3]), for exponential service times and general stationary arrival processes, which include, in particular, the interrupted Poisson process, Markov modulated Poisson Process (MMPP) and Markov arrival process (MAP). The main results are obtained by establishing multimodularity properties of the objective functions, which turns to be quite elegant and simpler than in other applications ([6, 2, 3]).

Using ideas from [10] we solve the dual problem [7] which is concerned with the sharing of a single communication link between multiple sources. This is modeled as an optimal assignment problem of a single server to several single buffer queues to which packets arrive according to Poisson processes. The relation to the previous problem is obtained if we take as states in the new problem the vector of free spaces in the buffers (instead of the occupancy of the buffers). The distribution of the process of free spaces in the new model is the same as the distribution of the process of buffer occupancy in the previous problem. This yields the solution of the service assignment problem.

The paper is structured as follows. In the next section we study the admission into a single server, and obtain the optimal policy. We then establish in the following Section properties of the cost and an ordering between the performance of different policies. Section 5 deals with the optimal routing problem. The service assignment problem is solved in section 6 and an application to a robot scheduling problem in the Web is discussed in Section 7. The technical arguments that establish the multimodularity properties are delayed to Section 8.

## 2 The admission into a single server

Consider a single server with no waiting room. Let $T_n$ be the point process representing the arrival epochs, and assume that service times are i.i.d., independent of the arrival process, and exponentially distributed with parameter $\mu$. We assume that $\tau_n := T_{n+1} - T_n$ is a stationary process (in $n$). An arrival can be rejected or accepted by an admission controller. An admission policy is a sequence $a = (a_1, a_2, ...)$ where $a_i = 1$ means acceptance of the $i$th arrival, and $a_i = 0$ means its rejection. The actions are taken without any knowledge of the state of the buffer, and if it is already full when the packet is admitted, then the packet is lost. Let $X_n = X_n(a)$ be the number of packets in the system just after the $n$th action is taken.

We consider first the following problem:
(P1) Maximize

$$g(a) \stackrel{\triangle}{=} \lim_{n \to \infty} \frac{1}{n} \sum_{j=1}^{n} E^a X_n$$

subject to the constraint that a fraction of no more than a fraction $p$ of the packets is accepted:

$$\lim_{n \to \infty} \frac{1}{n} \sum_{j=1}^{n} a_n \leq p. \tag{1}$$

The reason we consider this reward function for maximization is that maximizing the average number of packets in steady state is related to
- maximizing the throughput (the departure rate), and to
- minimizing the losses (due to both rejection by the controller and to the blocking).
Indeed, every customer has a sojourn time which is exponentially distributed with average $\mu^{-1}$, and therefore the average actual throughput is equal to the product of the long-run average number of packets in the system and $\mu$. The equivalence between maximizing the throughput and minimizing the losses follows since the loss rate is the initial given input rate (before the admission control and losses) minus the actual throughput.

Note that the queue size that we maximize in this section is that obtained by averaging over the times $T_n$. In some cases this will indeed coincide with the time average queue length (for example, if the interarrival times $\tau_n$ are exponentially distributed). In the case it does not coincide, we propose an alternative approach in Section 4.

Another motivation to study this problem follows from an interesting comparison with the infinite queue system. Below we show that for the current system the queue length is *maximized* by assignment sequences as regular as possible; for the infinite buffer system queue lengths are *minimized* by regular sequences ([3]). Some intuitive understanding can help to explain this at first sight contradictory phenomenon: in an infinite queue system minimizing queue lengths means minimizing waiting by spreading out arrivals; minimizing queue lengths in a system without buffers means maximizing losses by making arrivals bursty.

Let $e_i \in \mathbb{N}^m$, $i = 1, ..., m$ denote the vector having all entries zero except for a 1 in its $i$th entry. Define $d_i = e_{i-1} - e_i$, $i = 2, ..., m$. Let $\mathcal{F} = \{-e_1, d_2, ..., d_m, e_m\}$.

**Definition 2.1 (Hajek [6])** *A function $f$ on $\mathbf{Z}^m$ is multimodular with respect to $\mathcal{F}$ if for all $x \in \mathbf{Z}^m$, $v, w \in \mathcal{F}$, $v \neq w$,*

$$f(x + v) + f(x + w) \geq f(x) + f(x + v + w). \tag{2}$$

A sequence of $a_i, i = 1, ..., n$ is said to be feasible if all its components are either 0 or 1.

The *regular* of *balanced sequence* $\{a_k^p(\theta)\}$ with rate $p$ and initial phase $\theta$ is defined [6, 3] as,

$$a_k^p(\theta) = \lfloor kp + \theta \rfloor - \lfloor (k - 1)p + \theta \rfloor, \tag{3}$$

where $\lfloor x \rfloor$ is the largest integer smaller than or equal to $x$.

Before stating the main result, we present a simple coupling property:

**Lemma 2.1** *Fix an arbitrary policy $a$. Then one can construct a probability space $(\Omega, \mathcal{F}, P)$ such that*
*(i) the state trajectories starting from different initial states of the server couple in a time which is finite w.p. 1,*
*(ii) The following holds:*

$$\lim_{n \to \infty} |E[X_n(a)|X_0 = 1] - E[X_n(a)|X_0 = 0]| = 0.$$

4

**Proof.** We let the interarrival times be the same in both systems. If the policy $a$ never accepts packets or if arrivals never occur, then coupling occurs once the system starting at state 1 empties. Otherwise, Let $T_n$ be the time at which the first packet is admitted. Note that the admission occurs in both systems. If both systems are empty at time $T_n$, then coupling occurs at that instant (we take service times to be the same in both systems from time $T_n$ onwards). Otherwise, system 1 (in which the initial state is of a single packet) has 1 packet at time $T_n$, and system 0 (which is initially empty) has no packets. Coupling is obtained again at time $T_n$ by assuming that

(i) service times of all packets admitted after time $T_n$ are the same in both systems

(ii) At time $T_n$, the remaining service time of the packet in service in system 1 equals to the service time of the packet admitted in system 0 at time $T_n$.

Due to the memoryless property of the exponential distribution and the independence assumption on service times, the above coupling is consistent with the probability distribution of the original state processes. This establishes (i). (ii) follows from the bounded convergence theorem. ∎

**Theorem 2.1** *Assume that the system is controlled starting from time $T_1$. Assume that the interarrival times are stationary, and independent of the service times. Assume that the service times are i.i.d. exponentially distributed. Then for any $\theta \in [0,1]$, the balanced policy with rate $p$ and initial phase $\theta$ is optimal.*

**Proof.** Due to Lemma 2.1, we may assume without loss of generality that the system contains initially one packet at time $T_0$. $-E_{T,\sigma}X_n(a)$ is multimodular (we delay the proof of this property to Section 8). Here $\sigma$ denotes the random process governing the service completions. $-E_{T,\sigma}X_n(a)$ is clearly monotone decreasing in $a$. It remains to check conditions $<2>$ and $<3>$ in [1].

Condition $<3>$ requires for the multimodular sequence of functions $f_n(a)$ defined below that the following holds:

for any sequence $\{a_k\}$ $\exists$ a sequence $\{b_k\}$ such that

$\forall k, m, \quad k > m, \ f_k(b_1, \cdots, b_{k-m}, a_1, ..., a_m) = f_m(a_1, ..., a_m).$

This clearly holds in our case where $f_n(a) = -E_{T,\sigma}X_n(a)$ by setting $b_i = 1$. (The precise justification of the above is from property (i) which appears in Section 8.)

Condition $<2>$ states:

$f_k(a_1, ..., a_k) \geq f_{k-1}(a_2, ..., a_k), \ \forall k > 1;$

this holds in our case for $-E_{T,\sigma}X_n$ due to our assumptions on the initial state. ∎

# 3   Properties of the cost and of policies

**Lemma 3.1** $g(a) = g(\theta a)$ *for any policy* $a$*, where* $\theta$ *is the one step shift operator.*

**Proof.** This follows from Lemma 2.1. Indeed, with $a = (a_1, a_2, ...)$, we have

$$g(i) = \lim_{n \to \infty} \frac{1}{n} \sum_{j=1}^{n} E^a X_j = \lim_{n \to \infty} \frac{1}{n-1} \sum_{j=2}^{n} E^a X_j = \lim_{n \to \infty} \frac{1}{n} \sum_{j=1}^{n} E^{a'} X_j$$

where $a' = \theta a = (a_2, a_3, ....)$.  ∎

Define for any policy $a$

$$r_n(a) \triangleq \min\{m \geq 0 | a_{n-m} = 1\}.$$

**Lemma 3.2** *For any policy,* $EX_n(a)$ *is given by*

$$EX_n(a) = E \exp\left( -\mu \sum_{k=1}^{r_n(a)} \tau_k \right).$$

**Proof.** $X_n(a) = 1$ if and only if since time $T_{n-r_n(a)}$ there has been no departure. Thus $X_n(a) = 1$ if and only if during a time period of duration $\sum_{k=n-r_n(a)}^{n} \tau_k$ there are no end of services. Since $\tau_k$ is a stationary sequence, we obtain the wanted expression.  ∎

Assume that $a$ is periodic with period $S$. Let $a'$ be some shift of $a$ such that $a'_S = 1$. Thus $X_{nS}(a) = 1, n = 0, 1, 2, ....$ Then

$$g(a) = g(a') = \frac{1}{S} \sum_{k=1}^{S} E(X_k(a))$$

(see Sec. 7 in [4]).

Let $n = \sum_{k=1}^{S} a(k)$. Define

$$\eta(0) = 0; \qquad \eta(j+1) = \min\{l > \eta(j) : a_l = 1\}, j = 0, ..., n-1$$

Note that $\eta(n) = S$. Define

$$\delta_1 = \eta(1), \qquad \delta_j = \eta(j+1) - \eta(j), \ j = 2, ..., n-1.$$

6

**Lemma 3.3** *a (and thus a′) is fully determined by the sequence d up to a shift, and the cost g can be expressed as a function of $\delta(a)$:*

$$g(a) = \frac{1}{S} \sum_{i=1}^{n} \left( E \sum_{j=0}^{\delta_i - 1} e^{-\mu \sum_{k=0}^{j} \tau_k} \right) =: q(\delta) \tag{4}$$

**Proof.** By Lemma 3.2,

$$
\begin{aligned}
g(a) &= \frac{1}{S} \sum_{i=1}^{n} \sum_{j=0}^{\delta_i - 1} EX_{\eta(i)-j}(a) = \frac{1}{S} \sum_{i=1}^{n} \sum_{j=0}^{\delta_i - 1} E \exp\left( -\mu \sum_{k=0}^{r_{\eta(i)-j}(a)} \tau_k \right) \\
&= \frac{1}{S} \sum_{i=1}^{n} \sum_{j=0}^{\delta_i - 1} E \exp\left( -\mu \sum_{k=0}^{j} \tau_k \right) = \frac{1}{S} \sum_{i=1}^{n} \sum_{j=1}^{\delta_i} E \exp\left( -\mu \tau_1 \right)^j
\end{aligned}
$$

∎

We conclude that if $a$ is a periodic policy, then any policy obtained from $a$ by changing the order of the $\delta$ sequence that characterizes $a$, achieves the same cost. Thus a non-regular periodic policy of a period $(1, 0, 1, 0, 1, 0, 0, 1, 0, 0)$ has the same cost as a regular policy of the form $(1, 0, 1, 0, 0, 1, 0, 1, 0, 0)$.

This property can be seen to be a special corollary of a more general result that will be presented in Theorem 3.1 below.

Consider two $n$-dimensional vectors of integers $\delta(1), \delta(2)$.

**Definition 3.1** *(Majorization [9])*
*We say that $\delta(2)$ majorizes $\delta(1)$, which we denote by $\delta(1) \prec \delta(2)$, if*

$$
\begin{cases}
\displaystyle\sum_{i=1}^{k} \delta_{[i]}(1) \leq \sum_{i=1}^{k} \delta_{[i]}(2), & k = 1, ..., n-1, \\[2ex]
\displaystyle\sum_{i=1}^{n} \delta_{[i]}(1) = \sum_{i=1}^{n} \delta_{[i]}(2)
\end{cases}
$$

*where $\delta_{[i]}(j)$ is a permutation of $\delta_i(j)$ satisfying $\delta_{[1]}(j) \geq \delta_{[2]}(j) \geq ... \geq \delta_{[n]}(j)$, $j = 1, 2$.*
*A function g is called Schur convex (concave) if $g(\delta(1)) \leq (\geq)g(\delta(2))$ for each $\delta(1)$ and $\delta(2)$ such that $\delta(1) \prec \delta(2)$.*

**Lemma 3.4** *$q(\delta)$, defined in (4) is Schur concave.*

**Proof.** $q$ can be written as the sum

$$q(\delta) = \frac{1}{S} \sum_{i=1}^{n} \psi_i(\delta_i))$$

where $\psi$ is the term in brackets in eq. (4). Since

$$\psi_i(m+1) - \psi_i(m) = E e^{-\mu \sum_{k=0}^{m} \tau_k}$$

is monotone *decreasing* in $m$, it follows that $\psi_i$ are concave in $\delta_i$. The proof is then established by using proposition C.1 on p. 64 in [9]. ■

**Theorem 3.1** *Assume that $a$ and $a'$ are two periodic policies with the same period $S$ and the same sum*

$$n = \sum_{i=1}^{S} a_i = \sum_{i=1}^{S} a_i'.$$

*If $\delta(a') \prec \delta(a)$ then $g(a') \geq g(a)$.*

**Proof.** Follows from the Schur concavity of $g$, see [9] p. 54. ■

The above Theorem allows us to "improve" any given periodic policy by replacing it by a more "regular" one, where by more regular we mean a policy whose distance sequences are majorized by the less regular one.

A similar result was obtained in [5] in a related model (which we discuss in Section 7 under the assumption that the sequence $\tau_n$ is i.i.d.).

# 4   Time averages

In the previous sections we took averages of the costs as seen at times $T_n$, i.e. at arrival times. The problem with this is, however, that $T_n$ are in fact the times of *potential* rather than actual arrivals. In practice, arrivals occur only at a subsequence of $T_n$ which depend on the policy. In this section we obtain similar results for the actual arrival process.

We consider the same model of the system as well as the statistical assumptions as in Section 2.

Instead of describing a policy using a sequence $a$, it will be more helpful to consider an equivalent description using the distance sequence $\delta = (\delta_1, \delta_2, ...)$. Define $D(n) = \sum_{k=0}^{n} \delta_k$. The actual arrivals

occur at times $T_{D(n)}, n \in \mathbb{N}$. We define the process $\xi_n$ to be the number of packets in the buffer just prior to time $T_{D(n)}$. (If we took, as in the previous sections the time *after* a decision, then the number of packets would always be 1.)

The motivation for considering the system at arrival instants is the following. Whenever an actual arrival finds a packet in the system there is a loss. Thus minimizing the average number of packets at actual arrival times will also minimize the fraction of losses and maximize the throughput.

We consider first the following problem:
(Q1) Minimize

$$G(\delta) \triangleq \overline{\lim_{n \to \infty}} \frac{1}{n} \sum_{j=1}^{n} E^a \xi_n$$

subject to the constraint that a fraction of no more than a fraction $p$ of the packets is accepted, or in other words,

$$\overline{\lim_{n \to \infty}} \frac{1}{n} \sum_{j=1}^{n} \delta_n \geq 1/p. \tag{5}$$

Note that a policy $d$ is regular with rate $1/p$ if and only if its related policy $a$ is regular with rate $p$, see [2] Lemma 7.3.

**Theorem 4.1** *Assume that the system is controlled starting from time $T_1$. Assume that the inter-arrival times $\tau_n$ are stationary and independent of the service times. Assume that the service times are i.i.d. exponentially distributed. Then for any $\theta \in [0, 1]$, the balanced policy with rate $1/p$ and initial phase $\theta$ is optimal.*

The proof of the Theorem is similar to that of Theorem 2.1. The multimodularity of $E\xi_n(\delta)$ is established in Section 8.

We now present properties of the cost for periodic policies, which are the analogous of Section 3. We first note that the expected average cost for a periodic sequence $\delta = (\delta_1, ..., \delta_S)$ is given by

$$G(\delta) = \frac{1}{S} \sum_{i=1}^{S} E \exp\left(-\mu \delta_i\right) \tag{6}$$

when the interarrival times are i.i.d. Since this is a sum of functions that are convex in the $\delta_i$'s, it is Schur convex in $\delta$ (proposition C.1 on p. 64 in [9]). We thus conclude the following (see [9] p. 54):

9

**Theorem 4.2** *Assume that $\delta$ and $\delta'$ are two periodic policies with period $S$ and the same sum:*

$$n = \sum_{i=1}^{S} \delta_i = \sum_{i=1}^{S} \delta'_i.$$

*If $\delta' \prec \delta$ then $G(\delta') \prec G(\delta)$.*

# 5  Routing to several servers

We now consider $N$ servers (all with a single buffer), fed by a stationary input process as in Section 2. We make the same probabilistic assumptions as before on the service times in each server. Let $X_n^i$ be the number of packets (0 or 1) being served by server $i$ at the $n$th time epoch. A routing policy is a sequence $a = (a_1, a_2, ...)$ where $a_i = j$ means routing of the $i$th arrival to queue $j$. We consider the following routing problem:

(P2) Maximize

$$g(a) \triangleq \lim_{n \to \infty} \frac{1}{n} \sum_{j=1}^{n} \sum_{i=1}^{N} h_i E^a X_j^i,$$

where $h_i$ are some given positive constants.

The following theorems are the result of the properties we established in the previous section, the multimodularity properties (which is established in the next section) as well as the results in [1, 3].

**Theorem 5.1** *Consider the symmetric system, i.e., $\mu_1 = \cdots = \mu_N$. Then the round robin policy maximizes $g(a)$ (and hence, the expected average throughput).*

**Theorem 5.2** *Consider the case of two servers. Then there is some $p^*$ such that for any $\theta$, the policy that routes packets to server 1 according to the balanced policy with rate $p^*$ and initial phase $\theta$, and otherwise routes packets to server 2, is optimal.*

**Theorem 5.3** *Consider the case of two types of servers: a set $K_1 \subset \{1, ..., N\}$ of servers with $\mu = \mu_1$, and the remaining set $K_2$ of servers with $\mu = \mu_2$. Assume that $h_i = h^1$ are the same for all $i \in K_1$ and that $h_i = h^2$ for all $i \in K_2$. Then there exists some $p^*$ such that that for any $\theta$, the following policy is optimal:*
*it routes packets to the 1st group server according to the balanced policy with rate $p^*$ and initial*

*phase θ, and otherwise routes packets to the second group of servers. Within each group of servers, the order of service is round-robin.*

The proof of all three theorems is based on Proposition 2.22 in [3], which states that if a tuple $(p_1, p_2, ..., p_N)$ is made of less than (or of exactly) two distinct numbers, then it is balanceable. In other words, there exist a policy $a$ such that for each $i = 1, ..., N$, the sequence of indices in $a$ that correspond to routing arrivals to queue $i$ are is a regular with rate $p_i$. The optimality of balanceable sequences is established in Theorem 4.2 in [1]. (In fact, other cases are identified in Section 2 of [3] that yield balanceable sequences.)

**Remark 5.1** *For any $a$, let $\gamma^i(a) = \{\gamma_n^i(a), n \in \mathbb{N}\}$ be the sequence such that $\gamma_n^i(a) = 1$ if and only if $a_n = 1$. Using the result of Section 3, and in particular the Schur convexity of $g$ as a function of the δ's (see Lemma 3.4), one can show that if there are two periodic policies $a$ and $a'$ with the same period $S$, such that for any $i = 1, ..., N$,*

- $\sum_{j=1}^{S} \gamma_n^i(a) = \sum_{j=1}^{S} \gamma_n^i(a')$,

- $d(\gamma^i(a')) \prec d(\gamma^i(a))$,

*then $g(a') \geq g(a)$.*

# 6 The service assignment problem

Consider $N$ Poisson processes with parameters $\mu_1, ..., \mu_N$ respectively, of packet arrivals into $N$ respective single buffer queues. Only one buffer can obtain the transmission opportunity at a time. Let $T_n$ be time at which the $n$th transmission opportunity occurs. If an arrival occurs to a buffer that is full then it is lost; if the buffer is empty then the arriving packet is stored untill a transmission opportunity to that buffer arrives. If a buffer receives a transmission opportunity at time $T_n$ and it has a packet then this packet is transmitted, and immediately after time $T_n$, a new arriving packet can be stored in this buffer. If there is no packet in the buffer then this transmission opportunity is lost. We assume that $\tau_n := T_{n+1} - T_n$ is a stationary process (in $n$) and independent on the arrival process.

The role of the controller is to decide to which buffer the next transmission opportunity will be assigned. A service assignment policy is a sequence $a = (a_1, a_2, ...)$ where $a_i = j$ means that the

$i$th transmission opportunity will be to queue $j$. We assume that the controller has no information about the buffers' contents.

Let $Z_n^i$ be the number of packets (0 or 1) at buffer $i$ just after the $n$th action is taken. Denote $Y_n^i = 1 - Z_n^i$. It thus corresponds to the 'vacancies' process, as $Y_n^i$ equals one if the $i$th buffer is empty just after the $n$th action is taken (i.e. after time $T_n$). We consider the following problem: (P3) Minimize

$$g(a) \stackrel{\triangle}{=} \overline{\lim_{n \to \infty}} \frac{1}{n} \sum_{j=1}^{n} \sum_{i=1}^{N} h_i E^a Z_j^i,$$

where $h_i$ are some given positive constants.

The above objective corresponds to the minimization of blocking probabilities, since blocking occurs at queue $i$ between $T_n$ and $T_{n+1}$ if and only if $Z_n^i = 1$.

Note that by minimizing blocking probabilities, we maximize the throughput.

We now make the following key observation. Fix an arbitrary sequence $a$. Then the distribution of the vacancies process $\{Y_n^i\}_{n,i}$ in the service assignment problem is the same as the distribution of the buffer contents process $\{X_n^i\}_{n,i}$ in the routing problem. Hence, using the results of the previous section, we get the following main results.

**Theorem 6.1** *Consider the symmetric system, i.e., $\mu_1 = \cdots = \mu_N$. Then the round robin policy minimizes $g(a)$ (and hence, maximizes the expected average throughput).*

**Theorem 6.2** *Consider the case of two servers. Then there is some $p^*$ such that for any $\theta$, the policy that assigns transmission opportunities to buffer 1 according to the balanced policy with rate $p^*$ and initial phase $\theta$, and otherwise assigns transmission opportunities to server 2, is optimal.*

Similarly, one obtains the dual of Theorem 5.3.

# 7 Application to robot scheduling for Web search engines

We present in this section an application studied in [5] under assumptions more restrictive than ours.

The Web offers search engines, such as Alta Vista, Lycos, Infoseek and Yahoo, that serve as a data base that allow to search information on the WEB. These often contain robots that periodically

traverse the whole Web structure so as to update the data base. We consider the following problem: there is a fixed number $N$ of data-base pages. The contents of page $i$ are modified at times that follow a Poisson process with parameter $\mu_i$. A page is considered up-to-date by the Web engine from the time it is accessed by the robot untill the next time it is modified; at this point it is considered obsolete till the next time it is visited by the robot. The times between updates by the robot are given by a sequence $\tau_n$. In [5] these were assumed to be i.i.d.; in our framework we may allow them to form a general stationary ergodic sequence.

The problem is to find a visiting schedule for the robot that minimizes the obsolescence rate of the data base, i.e. the fraction of time during which the page $i$ is out of date.

The fact that the distribution of $\tau_n$ does not depend on the identity of the page that is updated (as is assumed in [5])) is not always the case, but it is justified in some applications. For example, if satellite links are used, then the main cause of delays is the earth to satellite distance, rather than the geographic location of the Web page that is accessed.

We now show that this problem is equivalent to the service assignment problem. Indeed, the robot can be considered as a single server that requires time $\tau_n$ in order to complete the $n$th update. The durations of consecutive modifications of page $i$ correspond to packet interarrival times to queue $i$. Finally, the time ratio of obsolescence of page $i$ corresponds in our service assignment problem to the blocking probability in queue $i$. (Indeed, blocking in queue $i$ is defined to be the period from the first time that a customer arrives to queue $i$ after a service there, till the next time that this queue is served.) Thus the results of the previous section holds for this model as well.

# 8   The multimodularity

**Lemma 8.1** *The two following statements are true. (i) $E_\sigma(-X_n(a))$ is multimodular,*
*(ii) $E_\sigma \xi_n(\delta)$ is multimodular.*

**Proof.** (i) The proof is based on the following useful properties. Let $q(a) \triangleq \max_{m \leq n}$ such that $a_m = 1$. Consider two policies $a, a'$. Then,

   Property (A): if $q(a) = q(a')$ then $X_n(a) =_{st} X_n(a')$ (this is equivalent to $E_\sigma(X_n(a)) = E_\sigma(X_n(a'))$ and to $P(X_n(a) = 1) = P(X_n(a') = 1)$).
This is due to the memoryless property of the exponential service times. Thus we can replace our original system by one where at each acceptance, the new packet replaces the one in service, instead of being rejected; the distribution of the process $X_n$ will not change.

<u>Property (B)</u>: if $q(a) \leq q(a')$ then $X_n(a) \leq_{st} X_n(a')$ (which is equivalent to $E_\sigma(X_n(a)) \leq E_\sigma(X_n(a')))$.

This is obtained by a similar argument.

We shall now check relation (2) for any $a$ such that $a, a+w, a+v, a+w+v$ are feasible.

Let $v = e_n$. We have for any $w \in \mathcal{F}$, $w \neq v$,

$$X_n(a+w+v) = X_n(a+v) = X_n(v) = 1,$$

and

$$X_n(a+w) \leq_{st} X_n(a).$$

The first relation follows from property (A) above, and the second from property (B), since

$$q(a+d_i) \leq q(a), i = 2, ..., n \quad \text{and} \quad q(a-e_i) \leq q(a)$$

(recall that $d_i$ corresponds to shifting an arrival to the past). This implies that

$$E_\sigma(-X_n(a+w+v)) = E_\sigma(-X_n(a+v)), \qquad E_\sigma(-X_n(a+w)) \geq E_\sigma(-X_n(a)).$$

Hence relation (2) is satisfied.

Next, assume $v = -e_1$. We consider $w \neq v$ (and thus restrict to $n \geq 2$). For $a+v$ to be feasible, we must have $a_1 = 1$. For $a+w$ to be feasible, we have $w \neq d_2$, and $q(a) > 1$. It then follows from property (A) that

$$X_n(a) =_{st} X_n(a+v), \qquad X_n(a+w) =_{st} X_n(a+v+w).$$

Hence relation (2) holds with equality.

It remains to check $v = d_i, d_j = w$, with $j > i$. Since $a + d_j$ is feasible, it follows that $q(a) \geq j$. Hence, by property (A),

$$X_n(a+d_i) =_{st} X_n(a).$$

Similarly, $q(a+d_j) \geq j-1$; since it is feasible then $a_{j-1} = 0$, so that $i < j-1$ (in order for $a + d_i$ to be feasible, we have to exclude $i = j-1$, since $a_{j-1} = 0$). Hence

$$X_n(a+d_i+d_j) =_{st} X_n(a+d_j).$$

Hence relation (2) holds with equality, which concludes the proof of (i).

(ii) Let $v$ be any one of the vectors in the set $\{-e_1, d_2, ..., d_{n-1}\}$. Then due to Property (A), for any $w \neq v$,

$$\xi_n(a) =_{st} \xi_n(a+v), \quad \xi_n(a+w) =_{st} \xi_n(a+v+w).$$

14

Hence relation (2) is satisfied. By symmetry it holds for any $w$ in the set $\{-e_1, d_2, ..., d_{n-1}\}$. It thus remains to check the case $v = e_n, w = d_n$.

From Lemma 3.2 we have:

$$E\xi_n(\delta) = E \exp\left(-\mu \sum_{k=1}^{\delta_n} \tau_k\right).$$

Let $x_n = m$, and let denote

$$y = \exp\left(-\mu \sum_{k=2}^{m} \tau_k\right).$$

Then

$$
\begin{aligned}
E\xi(x+v+w) = E\xi(x) &= Eye^{-\mu\tau_{m+1}} = Eye^{-\mu\tau_1} \\
E\xi(x+v) &= Eye^{-\mu[\tau_{m+1}+\tau_{m+2}]} = Eye^{-\mu[\tau_1+\tau_{m+1}]} \\
E\xi(x+w) &= Ey.
\end{aligned}
$$

Since the function $f(x) := ye^{-\mu x}$ is convex in $x$, it follows that $f(\tau_1 + z) - f(z)$ is increasing in $z$, so that

$$f(\tau_1 + \tau_{m+1}) - f(\tau_{m+1}) \geq f(\tau_1) - f(0).$$

By taking expectations, this implies relation (2) for $E\xi_n(a)$, which concludes the proof. ∎

# References

[1] E. Altman, B. Gaujal and A. Hordijk, "Multimodularity, Convexity and Optimization Properties", *Mathematics of Operations Research*, 25:324–347, 2000.

[2] E. Altman, B. Gaujal and A. Hordijk, "Admission Control in Stochastic Event Graphs", *IEEE Transactions on Automatic Control*, 45, 2000.

[3] E. Altman, B. Gaujal and A. Hordijk, "Balanced Sequences and Optimal Routing", *Journal of the ACM*, 2001 (to appear).

[4] E. Altman and A. Hordijk, "Applications of Borovkov's Renovation Theory to Non-Stationary Stochastic Recursive Sequences and their Control", *Advances of Applied Probability* **29**, pp. 388-413, 1997.

[5] E. G. Coffman Jr, Z. Liu, R. R. Weber, "Optimal robot scheduling for web search engines", *Journal of Scheduling*, 1:15–29, 1998.

[6] B. Hajek, "Extremal splitting of point processes", *Mathematics of Operations Research*, 10:543–556, 1985.

[7] A. Itai and Z. Rosberg, "A golden ratio control policy for a multiple-access channel", *IEEE Transactions on Automatic Control*, 29:712–718, 1984.

[8] G.M. Koole, "On the static assignment to parallel servers", *IEEE Transactions on Automatic Control*, 44:1588–1592, 1999.

[9] A. W. Marshall and I. Olkin, *Inequalities: Theory of of Majorization and Its Applications*, Academic Press, 1979.

[10] P.D. Sparaggis, C.G. Cassandras, and D.F. Towsley, "On the duality between routing and scheduling systems with finite buffer spaces", *IEEE Transactions on Automatic Control*, 38:1440–1446, 1993.